

Empirical Bayes

An example from practice

Disclaimers

1. I don't know much theory.
2. I will oversimplify / obfuscate many Google-specific details.
3. All data is (obviously) simulated.

Questions to Keep in Mind

What is Bayesian about Empirical Bayes?

When can you use it?

How can you test its assumptions?

Framework

Google search results for "mesothelioma". The search bar shows "mesothelioma" and the search button is a magnifying glass. Below the search bar are tabs for "Web", "Images", "Maps", "Shopping", "Videos", "More", and "Search tools". The results show "About 3,850,000 results (0.43 seconds)".

Ads related to **mesothelioma** ⓘ

- Mesothelioma Compensation - MesotheliomaClaimsCenter.info**
www.mesotheliomaclaimscenter.info/ 1 (866) 942 9878
Mesothelioma? Get Money You Deserve Fast! Get Help with Filing a Claim.
Mesothelioma Compensation Amounts - Help for Mesothelioma Victims
- Mesothelioma & Asbestos - Do you Have a Mesothelioma Lawsuit?**
www.lawfirm.com/Free_Consult 1 (855) 992 0102
Get the Settlement You Deserve.
Maximizing Mesothelioma Settlements - What To Expect Filing A Law Suit
- Mesothelioma Diagnosis? - We Help 50 People A Month**
www.veterans-mesothelioma.org/ 1 (888) 888 5043
New Information for US Veterans.
- Mesothelioma - Wikipedia, the free encyclopedia**
en.wikipedia.org/wiki/Mesothelioma ▼
Mesothelioma (or, more precisely, malignant mesothelioma) is a rare form of cancer that develops from cells of the mesothelium, the protective lining that covers ...
Asbestos - Mesotheliom - Peritoneal mesothelioma - Category:Mesothelioma

Ads ⓘ

- Mesothelioma Law Firm**
www.ghasites.com/ ▼
1 (888) 848 8991
Over 25 yrs experience. Billions in claims. Call us now or live chat.
- National Claims Center**
www.nationalmesotheliomaclaims.com/ ▼
1 (800) 713 6692
\$30B Asbestos Fund - BBB Approved
File Your Mesothelioma Claim Today
- Mesothelioma**
www.mesotheliomalawyer.co/ ▼
1 (888) 888 4056
Asbestos Claim Funds Available.
For California Residents Only.
- Mesothelioma**
www.mesothelioma-attorney-locators.com/ ▼
1 (800) 314 2433
Easily Find Mesothelioma Attorneys

Ads

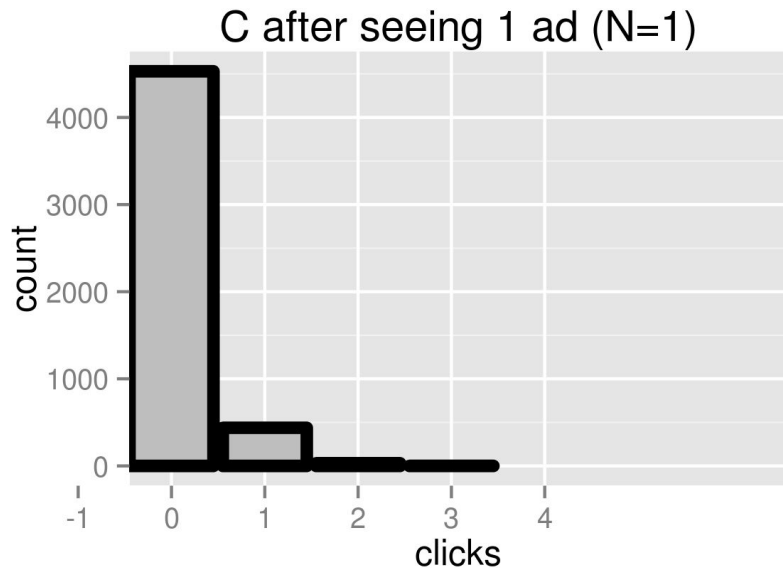
In a certain time period, an internet user sees N ads. Some they click on, some they don't.

For each user, count:

N = # of ads

C = # of clicks

A Dumb Question



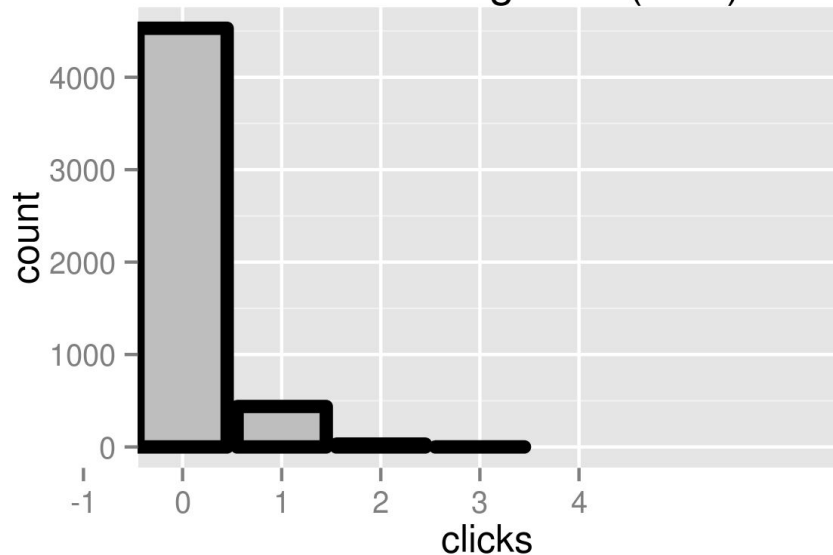
Look at 5000 people, each of whom saw one ad. The histogram of their clicks is shown.

90.6% had no clicks.

Does this mean 90.6% of people never click on ads?

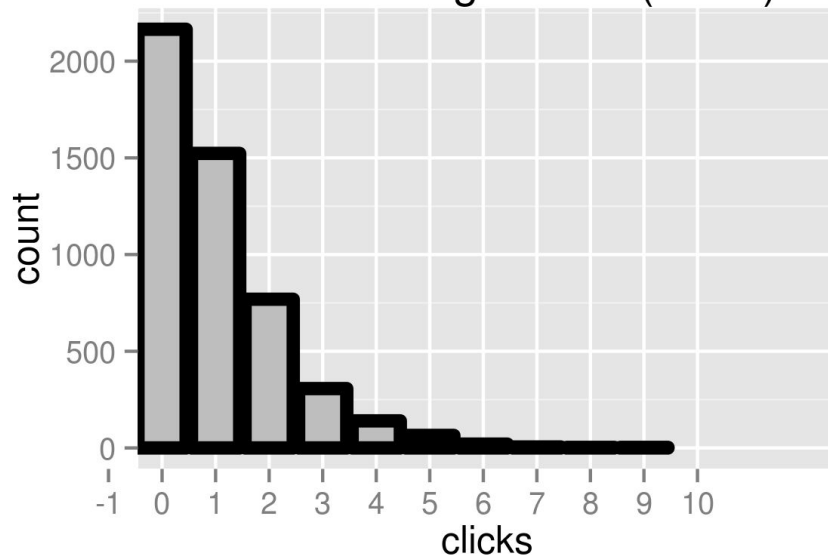
Same People, More Ads

C after seeing 1 ad (N=1)



90.6% no click

C after seeing 10 ads (N=10)



43.3% no click

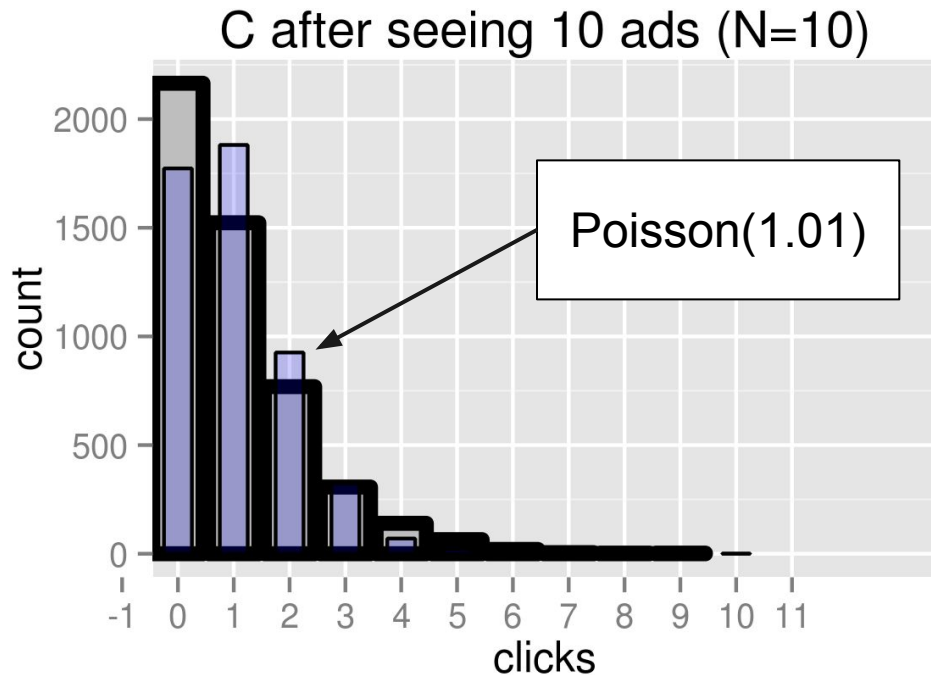
Clicks = Person + Noise

Clicks are a *noisy* observation of an attribute of a person at a particular time. Define:

p_i = the probability that person i clicks on an ad.

$$(C | p_i) \sim \text{Binom}(N, p_i) \approx \text{Poisson}(Np_i)$$

Overdispersion



Are all the p_i the same?
If so, because C is Poisson,
 $E(C) = \text{Var}(C) = Np_i$

But here,
C sample mean = 1.01
C sample var = 1.50

Mixture Model

Suppose the p_i are actually from a distribution.

$$p_i \sim \text{Gamma}(a, b)$$

$$\begin{aligned} C &\sim \text{Integral}(P(C|p_i) * P(p_i)) dp_i \\ &= \text{Negative Binomial}(c, d) \end{aligned}$$

Gamma Mixture of Poisson

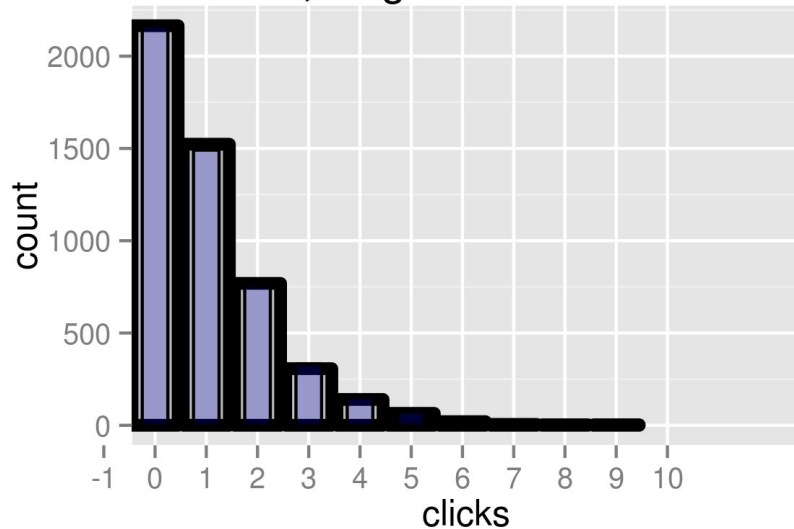
$$\begin{aligned} f(k; r, p) &= \int_0^\infty f_{\text{Poisson}(\lambda)}(k) \cdot f_{\text{Gamma}(r, \frac{p}{1-p})}(\lambda) \, d\lambda \\ &= \int_0^\infty \frac{\lambda^k}{k!} e^{-\lambda} \cdot \lambda^{r-1} \frac{e^{-\lambda(1-p)/p}}{\left(\frac{p}{1-p}\right)^r \Gamma(r)} \, d\lambda \\ &= \frac{(1-p)^r p^{-r}}{k! \Gamma(r)} \int_0^\infty \lambda^{r+k-1} e^{-\lambda/p} \, d\lambda \\ &= \frac{(1-p)^r p^{-r}}{k! \Gamma(r)} p^{r+k} \Gamma(r+k) \\ &= \frac{\Gamma(r+k)}{k! \Gamma(r)} p^k (1-p)^r. \end{aligned}$$

A gamma mixture of
poissons is a negative
binomial.

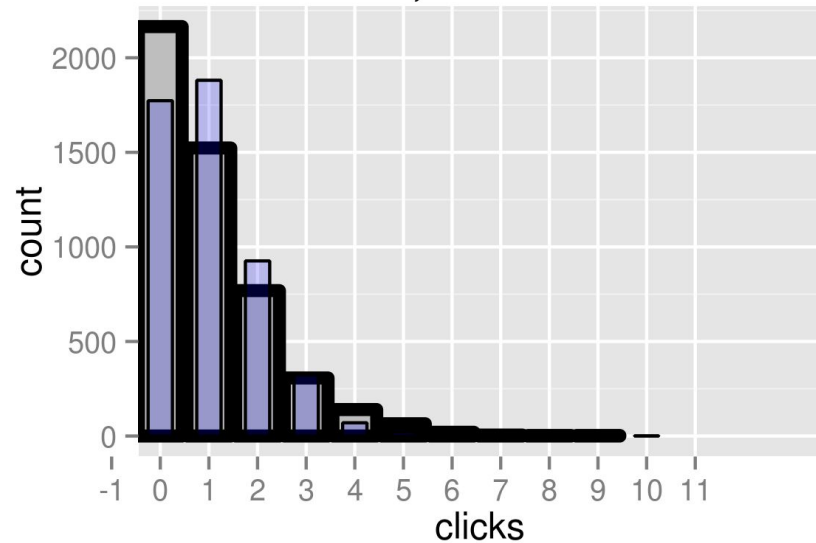
(Formula from Wikipedia)

Negative Binomial Fit

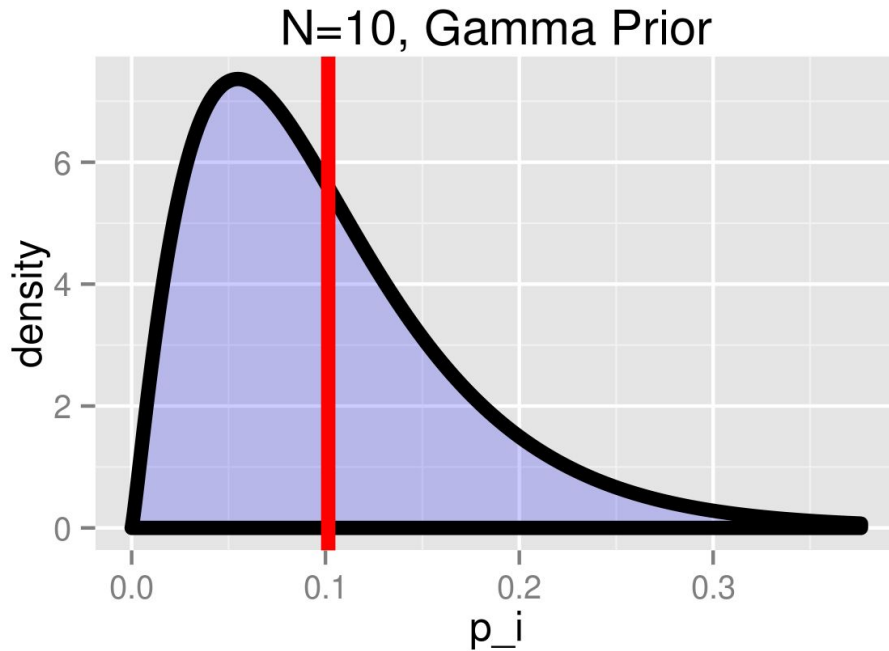
N=10, Negative Binomial Fit



N=10, Poisson Fit



Gamma Prior



Here is what the fit “prior” looks like.

This is an estimate of the distribution of click probabilities based on the overdispersion relative to Poisson.

Bayes' Rule

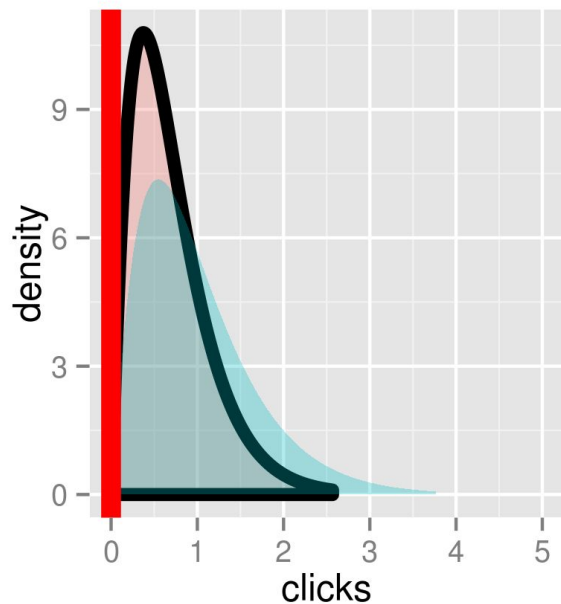
We draw a person from the urn and observe a C .

What might their p_i be?

- We know (have estimated) $P(p_i)$.
- We know (have assumed) $P(C | p_i)$.
- What we want is: $P(p_i | C) = P(C | p_i) P(p_i) / P(C)$

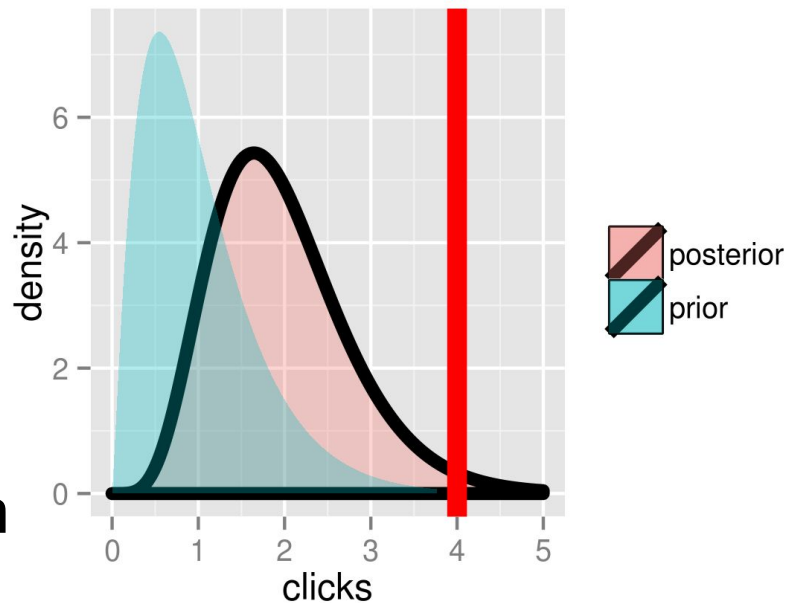
Lucky us! For gamma / poisson, $P(p_i | C)$ is a gamma distribution (gamma is conjugate for poisson).

Bayes' Rule



Belief about p_i when $C=0$

**Note:
regression
to the
mean**



Belief about p_i when $C=4$

Some Criticism - Meaning

- It is tempting to claim to have measured something about a person. Is this necessarily true? (hint: no)
- How can we test this?

Some Criticism - Meaning

OK	Day 1	Day 2
Person 1	C=9	C=12
Person 2	C=1	C=2

NOT OK	Day 1	Day 2
Person 1	C=9	C=2
Person 2	C=1	C=12

A key assumption was that a conditional on i , clicks are Poisson. We knew the noise, so we could detect overdispersion. If this is violated, EB will mislead you.

Some Criticism - Meaning

- Validate your conditional noise model
 - e.g. here, check that an *individual* is not over-dispersed
- Modeling explicitly can be hard
 - There is generally a nasty integral
- What will you do with it? What is your loss?
 - Sometimes you can measure this empirically.

Some Criticism - The Prior

- Ryan, your choice of prior seems suspiciously convenient.
- What if the fit is not so good?

Some Criticism - The Prior

For certain things, it is enough to fit the marginal:

- e.g. Robbins formula (?), cf stat 210 HW3:

$$h(z, \mu) = g(\mu)f(z; \mu) = g(\mu) \exp(z\mu - A(\mu))f_0(z)$$

$$\mathbb{E}(\mu|Z = z) = \frac{d}{dz} \left(\frac{f}{f_0} \right) (z) \bigg/ \left(\frac{f}{f_0} \right) (z) = \frac{d}{dz} \left(\log \left(\frac{f}{f_0} \right) \right) (z).$$

For some things, you are obviously out of luck.

- e.g. does $P(p_i=0.24572)$ exactly?

Some Criticism - The Prior

Other ways to get priors:

- Point mixtures of conjugate priors (EM)
- Method of moments
- Fit the marginal non-parametrically
 - Deconvolve
 - Use something like Robbins' formula
- Remember, you eventually need a posterior.

Questions Revisited:

What is Bayesian about Empirical Bayes?

When can you use it?

How can you test its assumptions?